

METHOD, SYSTEM, AND COMPUTER PROGRAM PRODUCT FOR
DETERMINING PROPERTIES OF COMBINATORIAL LIBRARY
PRODUCTS FROM FEATURES OF LIBRARY BUILDING BLOCKS

Inventors: Victor S. Lobanov
Dimitris K. Agrafiotis
F. Raymond Salemme

CROSS-REFERENCE TO RELATED APPLICATIONS

[0001] This application claims the benefit of U.S. Provisional Application No. 60/226,682, filed August 22, 2000, U.S. Provisional Application No. 60/235,937, filed September 28, 2000, and U.S. Provisional Application No. 60/274,238, filed March 9, 2001, each of which is incorporated by reference herein in its entirety.

[0002] The following application of common assignee is related to the present application, and is incorporated by reference herein in its entirety:

[0003] "System, Method and Computer Program Product For Fast and Efficient Searching of Large Chemical Libraries," serial number 09/506,741, filed February 18, 2000.

FIELD OF THE INVENTION

[0004] The present invention relates to combinatorial chemistry and computer aided molecular design. The present invention also relates to pattern analysis, information representation, information cartography and data mining. In particular, the present invention relates to predicting measurable or computed properties of products in a combinatorial chemical library based on features of their corresponding reagents.

BACKGROUND OF THE INVENTION

[0005] Algorithmic efficiency has been a long-standing objective in computational drug design. There is perhaps no other problem in chemistry where the need for efficiency is as pressing as in combinatorial chemistry. As will be understood by a person skilled in the relevant art, a significant

bottleneck in the virtual screening of a large combinatorial chemical library is the explicit enumeration of products and the calculation of their pertinent properties.

[0006] Whether it is based on molecular diversity, molecular similarity, structure-activity correlation, or structure-based design, the design of a combinatorial experiment typically involves the enumeration of every possible product in a virtual library, and the computation of key molecular properties that are thought to be pertinent to the application at hand. (See, e.g., Agrafiotis, D. K., The diversity of chemical libraries, *The Encyclopedia of Computational Chemistry*, Schleyer, P. v. R., Allinger, N. L., Clark, T., Gasteiger, J., Kollman, P. A., Schaefer III, H. F., and Schreiner, P. R., Eds., John Wiley & Sons, Chichester, 742-761 (1998); and Agrafiotis, D. K., Myslik, J. C., and Salemme, F. R., Advances in diversity profiling and combinatorial series design, *Mol. Diversity*, 4(1), 1-22 (1999), each of which is incorporated by reference herein in its entirety).

[0007] Several product-based methodologies for screening virtual libraries have been developed. (See, e.g., Sheridan, R.P., and Kearsley, S.K., Using a genetic algorithm to suggest combinatorial libraries, *J. Chem. Info. Comput. Sci.*, 35, 310-320 (1995); Weber, L., Wallbaum, S., Broger, C., and Gubernator, K., Optimization of the biological activity of combinatorial compound libraries by a genetic algorithm, *Angew. Chem. Int. Ed. Eng.*, 34, 2280-2282 (1995); Singh, J., Ator, M. A., Jaeger, E. P., Allen, M. P., Whipple, D. A., Solowej, J. E., Chowdhary, S., and Treasurywala, A. M., Application of genetic algorithms to combinatorial synthesis: a computational approach for lead identification and lead optimization, *J. Am. Chem. Soc.*, 118, 1669-1676 (1996); Agrafiotis, D. K., Stochastic algorithms for maximizing molecular diversity, *J. Chem. Info. Comput. Sci.*, 37, 841-851 (1997); Brown, R. D., and Martin, Y. C., Designing combinatorial library mixtures using genetic algorithms, *J. Med. Chem.*, 40, 2304-2313 (1997); Murray, C.W., Clark, D.E., Auton, T.R., Firth, M.A., Li, J., Sykes, R.A., Waszkowycz, B., Westhead, D.R. and Young, S.C., PRO_SELECT: combining structure-based drug design and combinatorial chemistry for rapid lead discovery. 1. Technology, *J. Comput.-Aided Mol. Des.*, 11, 193-207 (1997); Agrafiotis, D. K., and

Lobanov, V. S., An efficient implementation of distance-based diversity metrics based on k-d trees, *J. Chem. Inf. Comput. Sci.*, 39, 51-58 (1999); Gillett, V. J., Willett, P., Bradshaw, J., and Green, D. V. S., Selecting combinatorial libraries to optimize diversity and physical properties, *J. Chem. Info. Comput. Sci.*, 39, 169-177 (1999); Stanton, R. V., Mount, J., and Miller, J. L., Combinatorial library design: maximizing model-fitting compounds with matrix synthesis constraints, *J. Chem. Info. Comput. Sci.*, 40, 701-705 (2000); and Agraftotis, D. K., and Lobanov, V. S., Ultrafast algorithm for designing focused combinatorial arrays, *J. Chem. Info. Comput. Sci.*, 40, 1030-1038 (2000), each of which is incorporated by reference herein in its entirety).

[0008] These product-based methodologies become impractical, however, when they are applied to large combinatorial libraries, i.e. libraries that contain a large number of possible products. In such cases, the most common solution is to restrict attention to a smaller subset of products from the virtual library, or to consider each substitution site independently of all the others. (See, e.g., Martin, E. J., Blaney, J. M., Siani, M. A., Spellmeyer, D. C., Wong, A. K., and Moos, W. H., *J. Med. Chem.*, 38, 1431-1436 (1995); Martin, E. J., Spellmeyer, D. C., Critchlow, R. E. Jr., and Blaney, J. M., *Reviews in Computational Chemistry*, Vol. 10, Lipkowitz, K. B., and Boyd, D. B., Eds., VCH, Weinheim (1997); and Martin, E., and Wong, A., Sensitivity analysis and other improvements to tailored combinatorial library design, *J. Chem. Info. Comput. Sci.*, 40, 215-220 (2000), each of which is incorporated by reference herein in its entirety). Unfortunately, the latter approach, which is referred to as reagent-based design, often produces inferior results in terms of meeting the primary design objectives. (See, e.g., Gillet, V. J., Willett, P., and Bradshaw, J., *J. Chem. Inf. Comput. Sci.*; 37(4), 731-740 (1997); and Jamois, E. A., Hassan, M., and Waldman, M., Evaluation of reagent-based and product-based strategies in the design of combinatorial library subsets, *J. Chem. Inf. Comput. Sci.*, 40, 63-70 (2000), each of which is incorporated by reference herein in its entirety).

[0009] Hence there is a need for methods, systems, and computer program products that can be used to screen large combinatorial chemical libraries, which do not have the limitations discussed above.

SUMMARY OF THE INVENTION

[0010] The present invention provides a method, system, and computer program product for determining properties of combinatorial library products from features of library building blocks.

[0011] As described herein, at least one feature is determined for each building block of a combinatorial library having a plurality of products. A training subset of products is selected from the plurality of products of the combinatorial library, and at least one property is determined for each product of the training subset of products. A building block set is identified for each product of the training subset of products, and an input features vector is formed for each product of the training subset of products. A supervised machine learning approach is used to infer a mapping function that transforms the input features vector for each product of the training subset of products to the corresponding at least one property for each product of the training subset of products. After the mapping function is inferred, it is used for determining, estimating, or predicting properties of other products of the library. Properties of other products are determined, estimated, or predicted from their corresponding input features vectors using the inferred mapping function. Building block sets are identified for a plurality of additional products of the combinatorial library. Input features vectors are formed for the plurality of additional products. The input features vectors for the plurality of additional products are transformed using the mapping function to obtain at least one estimate property for each of the plurality of additional products.

[0012] In embodiments of the invention, both measured values and/or computed values are used as features for the building blocks of the combinatorial library. Both measured values and/or computed values are also used as properties for the products of the training subset. In embodiments of the invention, at least one of the features of the building blocks is the same as at least one of the properties of the products.

[0013] In an embodiment of the invention, the mapping function is implemented using a multilayer perceptron. The multilayer perceptron is trained to implement the mapping function using the input features vector and

the corresponding at least one property for each product of the training subset of products.

[0014] In an embodiment of the invention, the building blocks of the combinatorial library include reagents used to construct the combinatorial library. In other embodiments, the building blocks of the combinatorial library include fragments of the reagents used to construct the combinatorial library. In still other embodiments, the building blocks of the combinatorial library include modified fragments of the reagents used to construct the combinatorial library.

[0015] Further embodiments, features, and advantages of the present invention, as well as the structure and operation of the various embodiments of the present invention, are described in detail below with reference to the accompanying figures.

BRIEF DESCRIPTION OF THE DRAWINGS/FIGURES

[0016] The file of this patent contains at least one drawing executed in color. Copies of this patent with color drawing(s) will be provided by the Patent and Trademark Office upon request and payment of the necessary fee.

[0017] The present invention is described with reference to the accompanying drawings wherein:

[0018] FIGs.1A-B illustrate a flowchart of a method for determining properties of combinatorial products from a combinatorial library according to an embodiment of the present invention;

[0019] FIG.2 illustrates an example combinatorial neural network according to an embodiment of the present invention;

[0020] FIG. 3 illustrates average similarity scores for 10 sets of 1000 compounds most similar to each of 10 randomly chosen 'leads' as selected by various methodologies from a 6.29 million-member Ugi library according to the invention;

[0021] FIG. 4A illustrates a two-dimensional nonlinear map of compounds selected based on maximum similarity to a randomly chosen 'lead' using calculated product properties;

- [0022] FIG. 4B illustrates a two-dimensional nonlinear map of compounds selected based on maximum similarity to a randomly chosen 'lead' using estimated product properties according to the invention;
- [0023] FIG. 4C illustrates a magnified view of the area outlined in FIG. 4A;
- [0024] FIG. 4D illustrates a magnified view of the area outlined in FIG. 4B;
- [0025] FIG. 5 illustrates average similarity scores for 10 sets of 1000 compounds most similar to each of 10 randomly chosen 'leads' as selected by various methodologies from a 6.75 million-member diamine library according to the invention;
- [0026] FIG. 6 illustrates a comparison of central processing unit times required for similarity searching using a conventional methodology and combinatorial neural network methodologies according to the invention;
- [0027] FIG. 7 illustrates the quality of estimated product properties according to the invention comparing training and test sets of products selected from two combinatorial libraries;
- [0028] FIG. 8 illustrates the quality of estimated product properties according to the invention as a function of the training set size;
- [0029] FIGs. 9A-9E illustrate a table of example estimation of descriptor properties of combinatorial products by trained single-output networks according to the invention;
- [0030] FIG. 10 illustrates a table of average similarity scores and percent identity for 10 sets of 1000 compounds most similar to 10 randomly chosen 'leads' as selected by various methodologies from a 6.29 million-member Ugi virtual library according to the invention;
- [0031] FIG. 11 illustrates a reaction scheme for a 4-component combinatorial library based on the Ugi reaction;
- [0032] FIG. 12 illustrates a reaction scheme for a 3-component combinatorial library based on a two-step reductive amination reaction involving a diamine core and two sets of alkylating/acylating agents; and
- [0033] FIG. 13 illustrates an exemplary computing environment within which the invention can operate.

DETAILED DESCRIPTION OF THE INVENTION

[0034] Embodiments of the present invention are now described with references to the figures, where like reference numbers indicate identical or functionally similar elements. Also in the figures, the left most digit(s) of each reference number corresponds to the figure in which the reference number is first used. While specific configurations and arrangements are discussed, it should be understood that this is done for illustrative purposes only. One skilled in the relevant art will recognize that other configurations and arrangements can be used without departing from the spirit and scope of the invention. It will also be apparent to one skilled in the relevant art(s) that this invention can also be employed in a variety of other devices and applications.

Overview of the Invention

[0035] The present invention learns to determine, estimate, or predict values associated with properties of combinatorial library products from features of library building blocks. In operation, at least one feature of the building blocks is determined, retrieved, or obtained. A training subset of products is selected from the products, and values for at least one property is determined, retrieved, or obtained for each product of the training subset. A building block set is identified for each product of the training subset, and an input features vector is formed for each product of the training subset. A supervised machine learning approach is used to infer a mapping function that transforms the input features vector for each product of the training subset to the corresponding value of at least one property for each product of the training subset. After the mapping function is inferred, it is used for determining, estimating, or predicting properties of other products of the library from their corresponding input features vectors.

Method Embodiment of the Invention

[0036] In an embodiment, the present invention is applied to an electronic library of chemical compounds. The invention is not, however, limited to this example.

[0037] A combinatorial chemical library is a collection of chemical compounds or "products" generated by combining a number of chemical "building blocks" such as reagents. For example, a linear combinatorial chemical library such as a polypeptide library is formed by combining a set of chemical building blocks called amino acids in every possible or nearly every possible way for a given compound length (i.e., the number of amino acids in a polypeptide compound). Millions of products theoretically can be synthesized through such combinatorial mixing of building blocks. One commentator has observed that the systematic, combinatorial mixing of 100 interchangeable chemical building blocks results in the theoretical synthesis of 100 million tetrameric compounds or 10 billion pentameric compounds (Gallop *et al.*, "Applications of Combinatorial Technologies to Drug Discovery, Background and Peptide Combinatorial Libraries," J. Med. Chem. 37, 1233-1250 (1994), which is incorporated by reference herein in its entirety). As will be understood by a person skilled in the relevant art, a combinatorial library can be mathematically represented as combinatorial library P , $\{p_{l_1 l_2 \dots l_j}, i = 1, 2, \dots, r; j = 1, 2, \dots, r_i\}$, wherein r represents the number of variation sites in the combinatorial library P , and r_i represents the number of building blocks at the i -th variation site.

[0038] As used herein, the term "building blocks" refers to reagents, fragments of reagents, and/or modified fragments of reagents. In an embodiment of the invention, the building blocks of the combinatorial library comprise the reagents used to construct the combinatorial library. In other embodiments, the building blocks may comprise fragments of the reagents used to construct the combinatorial library and/or modified fragments of the reagents used to construct the combinatorial library.

[0039] FIGs. 1A and 1B illustrate a flowchart of the steps of a method 100 for determining, estimating, or predicting measurable or computable properties of

products in a combinatorial chemical library based on features of their corresponding reagents. Method 100 will now be described with reference to the steps illustrated in FIGs. 1A and 1B.

[0040] In step 110, at least one feature (descriptor) is determined for each building block of a combinatorial library having a plurality of products $\{a_{ijk}, i = 1, 2, \dots, r; j = 1, 2, \dots, r_i; k = 1, 2, \dots, n_i\}$, wherein r represents the number of variation sites in the combinatorial library P , r_i represents the number of building blocks at the i -th variation site, and n_i represents the number of features used to characterize each building block at the i -th variation site. As used herein, a feature value can be determined, for example, by computing a value or by retrieving a previously calculated or measured value from a storage medium.

[0041] In an embodiment of the invention, topological descriptors are computed as building block features. In another embodiment of the invention, the principal components required to capture 99% of the total variance are computed from the topological descriptors calculated for the building blocks. Other example descriptors or features that can be determined include quantum mechanical properties, pharmacophoric properties, BCUT properties and/or other molecular properties. Still other descriptors or features that can be determined will be known to persons skilled in the relevant arts given the description of the invention herein.

[0042] In an embodiment of the invention, at least one of the features of the building blocks is a calculated value. In another embodiment, at least one of the features of the building blocks is a measured value. In either embodiment, the feature values can be obtained or retrieved, for example, from an information storage device.

[0043] In step 120, a training subset of products is selected from the plurality of products of the combinatorial library. In an embodiment, a training subset of products $\{p_i, i = 1, 2, \dots, m; p_i \in P\}$ may be selected from a combinatorial library P . This training subset of products can be chosen in several manners. For example, the training subset of products can be chosen randomly. In another embodiment, the training subset of products can be chosen using a combinatorial design method. In yet another embodiment, the training subset

of products can be chosen using a diversity based selection technique. In a case of random selection, the composition of a particular training subset has little influence on the quality of an inferred mapping as long as the training subset is sufficiently large. As will be understood by a person skilled in the relevant arts given the description herein, the size of a training subset depends on the size of the combinatorial library and on the number of variation sites in the library in question.

[0044] In step 130, at least one property (descriptor) is determined for each product of the training subset of products. As used herein, a property value can be determined, for example, by computing or by retrieving a previously calculated or measured value from a storage medium. In an embodiment, q properties are determined for each compound p_i in the selected training subset of products, $y_i = \{y_{ij}, i = 1, 2, \dots, m, j = 1, 2, \dots, q\}$, wherein q is greater or equal to one.

[0045] In an embodiment of the invention, at least one of the properties of the products is a calculated value. In another embodiment, at least one of the properties of the products is a measured value. In either embodiment, the property values can be obtained or retrieved, for example, from an information storage device. In an embodiment, at least one of the features of the building blocks determined in step 110 is the same as at least one of the properties of the products determined in step 130. In another embodiment, none of the features of the building blocks determined in step 110 is the same as any of the properties of the products determined in step 130.

[0046] In step 140, a building block set is identified for each product of the training subset of products. As used herein, the term "building block set" refers to the at least one reagent, fragment of a reagent, and/or modified fragment of a reagent used to generate a product. The build block set for a particular product is referred to herein as corresponding to the product.

[0047] In an embodiment, the corresponding building blocks $\{t_{ij}, t_{ij} = 1, 2, \dots, r_j, j = 1, 2, \dots, r\}$ are identified for each product p_i of the training subset of products selected from the combinatorial library P .

[0048] In step 150, an input features vector is formed for each product of the training subset of products. As used herein, the term "input features vector"

refers to a single vector for a particular product of the combinatorial library formed by concatenating the features determined in step 110 for each of the one or more building blocks that make up the product's building block set. In an embodiment, building block features (e.g., reagent descriptors) are concatenated into a single array and presented to a combinatorial neural network according to the invention in the same order.

[0049] In an embodiment, for the combinatorial library P described above, input features vectors are formed by concatenating the features, determined in step 110, for the building blocks $\{t_{ij}, t_{ij} = 1, 2, \dots, r_j, j = 1, 2, \dots, r\}$ that are identified for each product p_i into a single vector $\{x_i = a_{1t_{i1}} | a_{2t_{i2}} | \dots | a_{rt_{ir}}\}$.

[0050] In step 160, a supervised machine learning approach is used to infer a mapping function that transforms the input features vector for each product of the training subset of products to the corresponding at least one property for each product of the training subset of products. In an embodiment, step 160 comprises the step of training a combinatorial neural network or a multilayer perceptron according to the invention, using the input features vector and the corresponding at least one property for each product of the training subset of products, to implement the mapping function. This may be represented mathematically as using a supervised machine learning approach to infer a mapping function f that transforms input values x_i to output values y_i from input/output pairs in a training set $T = \{(x_i, y_i), i = 1, 2, \dots, m\}$.

[0051] As described herein, embodiments of the invention uses a special class of neural networks, referred to herein as combinatorial networks or combinatorial neural networks (CNNs), that are trained to determine, estimate, or predict the properties of combinatorial products from the features of their respective building blocks. Generally speaking, a combinatorial network comprises an input layer containing $n_1 \times n_2 \times \dots \times n_r$ neurons, where r is the number of variation sites in the combinatorial library and n_i is the number of features used to characterize each building block at the i -th variation site. A typical combinatorial network may comprise one or more hidden layers that contain at least 2 neurons, depending on the complexity of the transformation,

and an output layer having a single neuron for each product feature predicted by the network.

[0052] In an embodiment of the invention, three-layer, fully connected multilayer perceptrons (MLPs) are used to form a combinatorial network. These networks can be trained using a standard error back-propagation algorithm (see, e.g., S. Haykin, Neural Networks, Macmillan, New York (1994), which is incorporated by reference herein in its entirety), and a logistic transfer function $f(x) = 1/(1 + e^{-x})$ can be used for both hidden and output layers. In accordance with the invention, each combinatorial network can be trained for a fixed number of epochs or until a predefined error threshold is met using, for example, a linearly decreasing learning rate from 1.0 to 0.01 and a fixed momentum of 0.8. During each epoch, the training patterns or samples can be presented to the network in a randomized order. In other embodiments, other combinatorial networks are used.

[0053] After a combinatorial network according to the invention is trained, analyzing or screening the combinatorial library (or any subset thereof) involves computing or retrieving precomputed features of building blocks, concatenating them into an input feature vector and feeding the input feature vector through the trained combinatorial network, which outputs estimate or predicted properties for the products. The estimate or predicted properties can then be used for any subsequent analysis, searching, or classification. As will be understood by a person skilled in the relevant art given the description herein, the present invention can be applied to a wide variety of molecular properties, regardless of origin and complexity.

[0054] Step 160 ends when the mapping function is inferred or a CNN is trained to implement the mapping function.

[0055] In step 170, building block sets are identified for a plurality of additional products of the combinatorial library. This step is similar to step 140 above.

[0056] In step 180, input features vectors are formed for the plurality of additional products. This step is similar to step 150 above.

[0057] In an embodiment, steps 170 and 180 involve identifying, after the mapping function f is determined, for a product $p_z \in P$, the corresponding reagents $\{t_{zj}, j = 1, 2, \dots, r\}$ and concatenating their features, $a_{1t_{z1}}, a_{2t_{z2}}, \dots, a_{rt_{zr}}$, into a single vector $\{x_z = a_{1t_{z1}} | a_{2t_{z2}} | \dots | a_{rt_{zr}}\}$.

[0058] In step 190, the input features vectors for the plurality of additional products are transformed using the mapping function of step 160 to obtain at least one estimate property for each of the plurality of additional products. This can be represented mathematically as mapping $x_z \rightarrow y_z$, using the mapping function (e.g., mapping function f) determined in step 160, wherein y_z represents the properties of product p_z . In embodiments of the invention, the estimate or predicted properties are stored for subsequent retrieval and analysis.

[0059] As will be understood by a person skilled in the relevant art given the description herein, in embodiments, the invention can be used to estimate or predict quantum mechanical properties of combinatorial compounds from quantum mechanical and/or other molecular properties of their respective building blocks. For example, the following quantum mechanical properties can be predicted according to the invention: molecular orbital energies; total electronic energy; total energy; heat of formation; ionization potential; and dipole moment.

[0060] In other embodiments, the invention can be used to predict BCUT properties (eigenvalues) of combinatorial compounds from BCUT and/or other molecular properties of their respective building blocks. As would be known to a person skilled in the relevant art, a BCUT value is an eigenvalue. As explained by R. S. Pearlman of the University of Texas, College of Pharmacy, the strength of intermolecular interactions depends on atomic charges, atomic polarizabilities, and atomic H-bond-abilities. Thus, Pearlman proposes constructing three classes of matrices to represent compounds: one class with atomic charge-related values on the diagonal, a second class with atomic polarizability-related values on the diagonal, and a third class with H-bond-abilities on the diagonal. Pearlman also proposed using a variety of additional definitions for the off-diagonal elements including functions of interatomic

distance, overlaps, computed bond-orders, etc. (See, e.g., R. S. Pearlman, *Novel Software Tools for Addressing Chemical Diversity*, at <http://www.netsci.org/Science/Combichem/feature08.html>.) According to Pearlman, the lowest and highest eigenvalues (i.e., BCUT values) of these matrices reflect aspects of molecular structure.

[0061] In embodiments, the invention can also be used to predict pharmacophoric properties of combinatorial compounds from pharmacophoric and/or other molecular properties of their respective building blocks. As would be known to a person skilled in the relevant art, a pharmacophore is the spatial mutual orientation of atoms or groups of atoms assumed to be recognized by and interact with a receptor or the active site of a receptor. A receptor can be envisioned as a macromolecular structure such as a protein, an enzyme or a polynucleotide being an integral part of the complex molecular structure of the cellular membrane in which it is anchored or associated with. The recognition elements or receptor sites are oriented in such a way that recognition of and interaction with ligands can take place, leading to a pharmacological effect.

[0062] As will be understood by a person skilled in the relevant art given the description herein, the invention is not limited to being used to predict just the above properties of combinatorial compounds from the properties of their respective building blocks. For example, the invention can be used to estimate or predict the 117 topological descriptors listed in FIGs. 9A-E. The invention can also be used to predict many other properties of combinatorial compounds from the properties of their respective building blocks.

Results and Discussion

[0063] In this section, the results obtained for embodiments of the method of the invention are presented and discussed. Three different combinatorial network architectures according to the invention were examined using the two combinatorial libraries described below. The network architectures examined were: (1) networks that take as input a single feature (descriptor) from each reagent and produce a single property (descriptor) for the product, (2)

networks that take as input multiple features (descriptors) from each reagent and produce a single property (descriptor) for the product, and (3) networks that take as input multiple features (principal components) from each reagent and produce a single property (principle component) for the product. The first architecture category is referred to herein as single-input single-output (SISO) perceptrons. The second and third architecture categories are referred to herein as multiple-input single-output (MISO) perceptrons.

[0064] The performance of each architecture was evaluated using three statistical measures: (1) the correlation coefficient between the actual and predicted properties of the products (descriptors), (2) the amount of distortion of the similarity matrix as measured by Pearson correlation coefficient, and (3) the effect of that distortion on similarity searching and context-based retrieval. As a person skilled in the relevant art would know, similarity searching represents the most common form of virtual screening. It is based on the 'similar property principle', i.e. the fundamental belief that structurally similar compounds tend to exhibit similar physicochemical and biological properties. (See Johnson, M. A., and Maggiora, G. M., *Concepts and Applications of Molecular Similarity*, Wiley (1990), which is incorporated by reference herein in its entirety). Thus, given a set of compounds with some desired biological effect, one seeks to identify similar compounds, expecting that some of them will be more potent, more selective, or more suitable in some other way than the original leads. For purposes of the evaluation described herein, the similarity between two compounds or products was measured by their Euclidean distance in the multidimensional space (see Willett, P.; Barnard, J.M.; Downs, G.M. Chemical Similarity Searching, *J. Chem. Info. Comput. Sci.*, 38, 983-996 (1998), which is incorporated by reference herein in its entirety) formed by the principal components that preserved 99% of the variance in the original topological features.

[0065] The simplest of the architectures involves a series of networks, each of which is trained to predict the value of a single product descriptor from the values of that descriptor of the corresponding reagents. Thus, for a library with r components, each product descriptor is estimated by a SISO network with r input and 1 output nodes, hereafter denoted $r-h-1$, where h is the number of

hidden nodes. This approach offers simplicity and ease of training, as well as access to the individual product descriptors. As illustrated in FIGs. 9A-9E, this embodiment of the invention works well for about 80% of the 117 topological descriptors used to evaluate the invention. About 20% of the descriptors listed in FIGs. 9A-9E were not predicted reliably using this embodiment.

[0066] The ability of CNNs according to the invention to estimate individual descriptors can be improved by increasing the number of synaptic parameters and by adding to the training data other reagent descriptors that can provide additional information needed for successful prediction. This leads to a network topology of the form $r \times n - h - 1$, where n is the number of input descriptors per reagent. The additional descriptors used with this embodiment of the invention can be chosen in a variety of ways. For example, one can employ a feature selection algorithm similar to that used in stepwise regression analysis. This involves trying all possible pairs of descriptors and select the best pair, then trying all possible triplets keeping the first two descriptors fixed and select the best triplet, and continuing in this manner until a predefined number of descriptors or error threshold is met. In practice, however, this rather intensive algorithm is unnecessary. Good results can be obtained using the following approach. First, the correlation coefficients between each reagent and each product descriptor are calculated, and a series of SISO networks are trained in the manner described herein. Then, for each product descriptor that cannot be adequately modeled (e.g., one having a training R^2 less than 0.9), the two reagent descriptors that are most highly correlated to that product descriptor are added to the training data, and a new MISO network is trained. When applied to the Ugi library (see FIG 11), this approach resulted in an array of neural networks that were able to predict all 117 descriptors with high accuracy for both the training and test sets (see FIGs. 9A-9E). As illustrated in FIGs. 9A-9E, the correlation coefficients between the actual and predicted descriptors ranged from about 0.77 to 1.0. The smaller values are typically associated with more complex properties such as the Bonchev-Trinajstic information index \overline{I}_D^C (see Bonchev, D. and

Trinajstić, N., *J. Chem. Phys.* 67, 4517-4533 (1977), which is incorporated by reference herein in its entirety) and the Kappa shape index $^3\chi_\alpha$ (see Hall L.H. and Kier, L.B, The Molecular Connectivity Chi Indexes and Kappa Shape Indexes in Structure-Property Relations, in *Reviews of Computational Chemistry*, Boyd, D.B. and Lipkowitz, K.B., Eds., VCH Publishers, Chapter 9, 367-422 (1991), which is incorporated by reference herein in its entirety).

[0067] To assess the impact on molecular similarity, the optimized networks were used in a feed-forward manner to estimate the descriptors of all 6.29 million compounds in the Ugi library. These descriptors were subsequently decorrelated using the rotation matrix derived from the training set, and the Pearson correlation coefficient of the resulting pairwise distances was computed. This statistic, which measures the correlation between the similarity coefficients computed with the two sets of descriptors (calculated vs. estimated), had a value of 0.99, indicating a nearly perfect reproduction. As used herein, the term 'calculated' refers to descriptors computed with the conventional method, and the term 'estimated' refers to the descriptors generated by the neural networks.

[0068] This accuracy was also reflected in the context of similarity searching, using 10 randomly chosen compounds from the Ugi library as 'leads.' In particular, the 1000 most similar compounds to each of these leads were identified using the PCs derived from both the calculated and estimated descriptors, and their similarity scores were compared. FIG. 10 shows a summary of the results obtained.

[0069] Note that in order to permit a direct comparison, the hit lists obtained with the estimated descriptors were fully enumerated, and their similarity scores were re-evaluated using calculated descriptors computed in the conventional manner. As shown in FIG. 3, in all 10 cases, the two designs had nearly identical scores and very similar content with an overlap ranging from 75 to 86 percent (see FIG. 10). The equivalence of these selections for one of the leads is graphically illustrated in the nonlinear maps of FIGs. 4A-4D. FIGs. 4A and 4C illustrate the case for the calculated descriptors. FIGs. 4B and 4D illustrate the case for the estimated descriptors. FIG.s 4C and 4D are magnified views of the areas outlined in FIGs. 4A and 4B.

[0070] The entire screening process, including enumeration of the training set, network training, decorrelation, and similarity searching, required only 35 minutes of CPU time. As illustrated in FIG. 6, this represents a 30-fold improvement in throughput compared to the direct approach.

[0071] Since principal components are often the desired output, significant improvements can be achieved if the evaluation of the individual descriptors are circumvented, and the combinatorial networks are trained to predict the principal components directly. As note herein, high-dimensional data sets are almost always redundant. For example, the 117 topological descriptors illustrated in FIGs. 9A-9E can be reduced to 25-30 latent variables without any significant loss in the contribution to variation. The presence of correlated variables affects molecular similarity in two important ways: (1) redundant features are effectively taken into account with a higher weight, and (2) there is a substantial and unnecessary increase in the computational effort required for data analysis.

[0072] The invention was evaluated using the combinatorial libraries described herein as follow. A sample set of 10,000 compounds was selected at random from the entire Ugi library, and was characterized using the set of 117 topological descriptors listed in FIGs. 9A-9E. These descriptors were normalized and decorrelated to 25 principal components, which accounted for 99% of the total variance in the data. In addition, all the reagents involved in making the entire Ugi library were described by the same set of descriptors, and were independently normalized and decorrelated to 27 principal components using the same variance cutoff.

[0073] These data were used to develop an array of 25 CNNs (denoted PC-MISO), each of which was trained to predict one of the product PCs using all 27 PCs from each of the 4 input reagents. Thus, each neural network was comprised of 108 input, 2 hidden, and 1 output neurons. Experiments showed that increasing the number of hidden neurons beyond two did not offer any significant improvements in the predictive ability of the resulting networks.

[0074] A set of 10,000 input-output pairs was randomly split into a training set containing 90% of the samples and a test set containing the remaining 10% of the samples. Each neural network was trained on the training set for 100

epochs or until a predefined error threshold was met. Once training was complete, the combinatorial networks were used in a feed-forward manner to predict the 25 PCs for all 6.29 million compounds in the Ugi library, which were, in turn, used to identify the 1000 most similar compounds to each of the 10 'leads' described herein.

[0075] The obtained selections were finally assessed using 'calculated' PCs and compared with the ideal solutions (see FIG. 10). Again, in all 10 cases, the selections were very similar to those derived with 'calculated' descriptors and slightly better than those derived with regular SISO and MISO CNNs, both in terms of their similarity scores and the identity of the selected compounds which ranged from 80-85% (see FIG. 10).

[0076] The entire screening process required only 39 minutes on an 800 MHz Pentium III processor.

[0077] In order to validate the generality of the invention, similar types of selections were carried out from a 3-component diamine library (see FIG 12), using the same set of 117 topological descriptors for both reagents and products. In this case, 29 and 28 PCs were necessary to capture 99% of the variance in the reagent and products descriptors, respectively. Thus, 3-3-1 SISO and 9-3-1 MISO networks were used to predict individual descriptors, and 87-3-1 PC-MISO networks were employed for the prediction of principal components.

[0078] As with the Ugi library, 10 leads were selected at random from the entire library and the 1000 most similar compounds to each of these leads were identified using the PCs derived from both the exact and approximate descriptors. Once again, the selections obtained with approximate PCs were virtually identical to the ideal solutions, with PC-MISO predictions leading to slightly better similarity scores (see FIG. 5).

[0079] The accurate reproduction of the similarity matrix is accompanied by an impressive gain in performance (see FIG. 6). Although for both libraries the training of SISO, MISO, and PC-MISO CNNs required comparable execution times, the latter performed slightly but consistently better. On the other hand, SISO and MISO networks provide access to individual descriptors, which may have additional utility in applications such as, for example,

diversity profiling, ADME modeling, and structure-activity correlation. Based on the evaluations described herein, networks with multiple output nodes (i.e. multiple-input multiple-output (MIMO) perceptrons producing multiple product descriptors or principal components) tend to be more difficult to train and produced results that are less accurate than those obtained with an ensemble of single-output networks.

[0080] As described above, a common concern with any machine learning algorithm is its dependence on the nature of the training set. To examine the effect of the composition of the training set on the quality of the predictions obtained by the CNNs described herein, 10 random samples of 10,000 compounds were drawn from the Ugi library and were used to train 10 different sets of 25 PC-MISO networks. The average R^2 between the pairwise distances computed with 'exact' and 'approximate' PCs over all 10 trials was 0.9951 ± 0.0004 and 0.9951 ± 0.0006 for the training and test set, respectively. The R^2 was computed by comparing the Euclidean distances between 1,000,000 randomly chosen pairs of compounds in the two PC spaces. Similar standard deviations were also observed for the diamine library (0.0003 and 0.0007 for the training and test set. (See FIG. 7.) This result suggests that the training of CNNs according to the present invention is both stable and convergent.

[0081] In a case of random selection, the size of the training set has a moderate effect on the quality of predictions as long as it remains large enough to sample each reagent sufficiently. The predictions improve as the size of the training set increases, and eventually plateaus after a few thousand samples (see FIG. 8). For the Ugi library there was virtually no improvement in prediction when the size of the training set was doubled from 10,000 to 20,000 compounds, but this was not the case for the diamine library where the difference in R^2 was still noticeable. The reason for this result is almost certainly related to the difference in the number of reagents involved in the construction of these libraries (254 for the Ugi and 400 for the diamine library) and the fact that, for a given sample size, each individual reagent is more extensively sampled in the Ugi library. The disadvantage of using larger training sets is that longer times are required for descriptor calculation and

network training. Thus, in general, the sample size should be determined by weighing the benefits of higher accuracy against the increasing cost of computation.

Combinatorial Libraries

[0082] In this section, the two example Ugi and diamine combinatorial libraries used to evaluate the invention are described.

[0083] The first combinatorial library used to evaluate the invention was a Ugi library containing 6.29 million compounds. FIG. 11 illustrates a reaction scheme 1100 for generating a 4-component combinatorial library based on the Ugi reaction. The Ugi library used to evaluate the invention was constructed using a set of 100 acids, 100 amines, 37 aldehydes, and 17 isonitriles chosen at random from the Available Chemicals Directory (MDL Information Systems, Inc., 140 Catalina Street, San Leandro, CA 94577).

[0084] The second combinatorial library used to evaluate the invention was a diamine library containing 6.75 million compounds. FIG. 12 illustrates a reaction scheme 1200 for a combinatorial library based on a two-step reductive amination reaction involving a diamine core and two sets of alkylating/acylating agents. The diamine library used to evaluate the invention was constructed using a set of 300 diamines and two sets of 150 alkylating/acylating agents selected at random from the Available Chemicals Directory.

[0085] The size of Ugi and diamine libraries was intentionally restricted so that an exhaustive search of these libraries would be possible in order to validate the results obtained using the method embodiment of the invention described herein. Both the reagents and the products of these libraries were characterized by a set of 117 topological descriptors, including molecular connectivity indices, kappa shape indices, subgraph counts, information-theoretic indices, Bonchev-Trinajstis indices, and topological state indices. These descriptors have a proven track record in structure-activity analysis, can be computed directly from the connection table, and are consistent with the medicinal chemists' perception of molecular similarity. Moreover, they have

been shown to exhibit proper 'neighborhood behavior' and are thus well suited for diversity analysis and similarity searching. These data were subsequently normalized and decorrelated using principal component analysis (PCA), resulting in an orthogonal set of 25 to 29 latent variables, which accounted for 99% of the total variance in the data. The PCA preprocessing step was necessary in order to eliminate duplication and redundancy in the data, which is typical of graph-theoretic descriptors.

[0086] For the nonlinear maps illustrated in FIGs. 4A-D, this multidimensional data was further reduced to two dimensions using the methodology described in U.S. Patent Application Ser. No. 09/823,977, filed April 3, 2001, titled "Method, System, And Computer Program Product For Representing Object Relationships In A Multidimensional Space," which is incorporated by reference herein in its entirety. The pair-wise distances between the points in the multidimensional principle component space are preserved on the two-dimensional nonlinear maps of FIGs. 4A-D. The two-dimensional nonlinear maps of FIGs. 4A-D were used to visualize the product selections described herein, which were carried out using all significant principle components.

Summary

[0087] As described above, the method of the invention can be used to estimate or predict properties of products using the features of reagents, thereby effectively eliminating the need to enumerate and describe every individual product in a virtual combinatorial chemical library. By circumventing enumeration and replacing descriptor evaluation with a simple feed-forward pass through a combinatorial neural network according to the invention, the invention permits the *in silico* characterization and screening of huge combinatorial libraries unmanageable by other means. Although the descriptors or properties produced by the invention are estimated values rather than calculated values, any differences between the estimated values of the invention and the calculated values obtained using conventional methods is minimal and has little or no impact on similarity searching. Embodiments of

the invention are more than an order of magnitude faster than conventional enumerative similarity searching methodologies, and this differential increases with the size and combinatorial complexity of the virtual library under investigation.

System and Computer Program Product Embodiments

[0088] As will be understood by a person skilled in the relevant arts given the description herein, the method embodiment of the invention described above can be implemented as a system and/or a computer program product. FIG. 13 shows an example computer system 1300 that supports implementation of the present invention. The present invention may be implemented using hardware, software, firmware, or a combination thereof. It may be implemented in a computer system or other processing system. The computer system 1300 includes one or more processors, such as processor 1304. The processor 1304 is connected to a communication infrastructure 1306 (e.g., a bus or network). Various software embodiments can be described in terms of this exemplary computer system. After reading this description, it will become apparent to a person skilled in the relevant art how to implement the invention using other computer systems and/or computer architectures.

[0089] Computer system 1300 also includes a main memory 1308, preferably random access memory (RAM), and may also include a secondary memory 1310. The secondary memory 1310 may include, for example, a hard disk drive 1312 and/or a removable storage drive 1314, representing a floppy disk drive, a magnetic tape drive, an optical disk drive, etc. The removable storage drive 1314 reads from and/or writes to a removable storage unit 1318 in a well-known manner. Removable storage unit 1318 represents a floppy disk, magnetic tape, optical disk, etc. As will be appreciated, the removable storage unit 1318 includes a computer usable storage medium having stored therein computer software and/or data. In an embodiment of the invention, removable storage unit 1318 can contain input data to be projected.

[0090] Secondary memory 1310 can also include other similar means for allowing computer programs or input data to be loaded into computer system

1300. Such means may include, for example, a removable storage unit 1322 and an interface 1320. Examples of such may include a program cartridge and cartridge interface (such as that found in video game devices), a removable memory chip (such as an EPROM, or PROM) and associated socket, and other removable storage units 1322 and interfaces 1320, which allow software and data to be transferred from the removable storage unit 1322 to computer system 1300.

[0091] Computer system 1300 may also include a communications interface 1324. Communications interface 1324 allows software and data to be transferred between computer system 1300 and external devices. Examples of communications interface 1324 may include a modem, a network interface (such as an Ethernet card), a communications port, a PCMCIA slot and card, etc. Software and data transferred via communications interface 1324 are in the form of signals 1328 which may be electronic, electromagnetic, optical or other signals capable of being received by communications interface 1324. These signals 1328 are provided to communications interface 1324 via a communications path (i.e., channel) 1326. This channel 1326 carries signals 1328 and may be implemented using wire or cable, fiber optics, a phone line, a cellular phone link, an RF link and other communications channels. In an embodiment of the invention, signals 1328 can include input data to be projected.

[0092] Computer programs (also called computer control logic) are stored in main memory 1308 and/or secondary memory 1310. Computer programs may also be received via communications interface 1324. Such computer programs, when executed, enable the computer system 1300 to perform the features of the present invention as discussed herein. In particular, the computer programs, when executed, enable the processor 1304 to perform the features of the present invention. Accordingly, such computer programs represent controllers of the computer system 1300.

Conclusion

[0093] While various embodiments of the present invention have been described above, it should be understood that they have been presented by way of example, and not limitation. It will be apparent to persons skilled in the relevant art that various changes in detail can be made therein without departing from the spirit and scope of the invention. Thus the present invention should not be limited by any of the above-described exemplary embodiments, but should be defined only in accordance with the following claims and their equivalents.

1503.1070003